

1                   **CONCURRENT DATA RECALL IN A HIERARCHICAL**  
2                   **STORAGE ENVIRONMENT USING PLURAL QUEUES**

3  
4                   **FIELD OF THE INVENTION**

5                   This invention relates generally to the field of computing and, more  
6 particularly, to a system and method for recalling data objects stored on media such as  
7 tapes or removable disks.

8  
9                   **BACKGROUND OF THE INVENTION**

10                  In computing systems that employ a mass storage device (such as a hard  
11 disk) for the storage of data objects (e.g., files), it is often the case that data objects  
12 stored on the mass storage device are "migrated" to backup media (e.g., tape, writeable  
13 optical disk, etc.) in order to make room for new data in mass storage. When a data  
14 object has been migrated to backup media, it may become necessary to restore the  
15 object from the backup media – for example, if a user requests read or write access to a  
16 migrated object. When access to such a "migrated" objects is requested, it may be  
17 necessary to locate the backup media containing the object, and to mount the media on  
18 an appropriate drive in order to restore the object.

19                  It is often the case that migrated data objects are distributed over several  
20 media (e.g., where data objects are files, files 1, 3, and 5 may be on tape 1, and files 2  
21 and 4 may be on tape 2). Conventional systems generally place restore requests in a  
22 single queue and process these requests in the order received (i.e., first in, first out)  
23 without regard to where the data is located. Processing restore requests in this manner  
24 tends to lead to repetitive mounting and dismounting of media, as well as repeated  
25 traversal of the same media. In the above example, if files 1, 2, 3, 4, and 5 are  
26 requested in that order and only one drive is available, then tapes 1 and 2 must  
27 alternately be mounted and dismounted from the drive after each file is restored.  
28 Moreover, even if two objects reside on the same medium, they may be requested (and

1 processed) in an order that differs from that in which they are located on the medium.  
2 In the case of a sequential medium such as a tape, this means that the tape may have to  
3 shuttle back and forth repeatedly in order to locate the requested items, thereby  
4 increasing wear and tear on the tape.

5 Additionally, it is usually the case that even if migrated objects reside on  
6 different media and plural drives are available, the objects are processed sequentially  
7 (i.e., one at a time) in the order in which they appear on a single queue. Restoring  
8 migrated objects in this manner is wasteful of resources, because one or more available  
9 drives may stand idle while data is retrieved from only a single drive.

10 The present invention overcomes the limitations and drawbacks of the  
11 prior art.

### 12 13 SUMMARY OF THE INVENTION

14 The invention provides a system and method for restoring data objects  
15 from backup media. Various objects to be restored may each reside on different media.  
16 When a migrated objects is requested, a database lookup is performed to determine on  
17 which medium the object is located. A queue is created for each medium that stores  
18 requested objects, and a request for a given object is placed in the queue corresponding  
19 to the object's host medium. Each queue may be either "active" or "non-active." A  
20 queue is "active" when its corresponding medium is mounted on a drive such that  
21 migrated files may be retrieved from that medium; otherwise, the queue is "non-  
22 active." An active queue is "processed" (i.e., the data objects on the queue are  
23 retrieved from the corresponding medium) until the queue is empty. Non-active queues  
24 wait until a drive becomes available, and are then processed in the same manner. New  
25 requests may be placed on both active and non-active queues.

26 The invention may be incorporated into a physical computing  
27 arrangement having more than one drive. In such a case, a number called the  
28 "concurrency" is defined, which is the number of drives that may be used concurrently

1 to read backup media. Any number of queues – up to the concurrency number – may be  
2 active at a given point in time. When the number of active queues equals the  
3 concurrency number, this means that the maximum allowable number of media are  
4 mounted in drives and are being used to restore migrated files. The concurrency  
5 number may be equal to the number of physical drives available, or it may be less than  
6 the number of physical drives (e.g., in the case where the system administrator wants to  
7 reserve one or more drives for other uses). Preferably, when plural queues are active at  
8 the same time, the corresponding media are concurrently read from different drives,  
9 thereby increasing the throughput of the restoration process by allowing different  
10 backup media to be read at the same time.

11 It is preferable that items be placed on the queues in an order based on  
12 where the requested items are located on the corresponding medium, in order to  
13 minimize traversal of the medium. For example, when sequential media such as tapes  
14 are used (or other media whose storage locations are traversed in a pre-defined logical  
15 sequence), the queues can be organized in monotonically increasing sequences based on  
16 the requested data object's offset relative to a start position. If a newly requested item is  
17 located ahead of the current position of the tape head (i.e., the new item has not yet  
18 been encountered as the tape moves in the forward direction), it is placed on the queue  
19 in the first sequence; if a newly requested item is behind the tape head (i.e., the items  
20 position has already been encountered as the tape moves forward), then it is placed in  
21 the second sequence. Once the first sequence is exhausted, the tape can be rewound so  
22 that the second sequence can be started, thereby reducing repeated back-and-forth  
23 traversal of, and therefore wear and tear on, the tape.

24 Other features of the invention are described below.

25

## 26 BRIEF DESCRIPTION OF THE DRAWINGS

27 The foregoing summary, as well as the following detailed description of  
28 preferred embodiments, is better understood when read in conjunction with the

Mon  
Snig

1 appended drawings. For the purpose of illustrating the invention, there is shown in the  
2 drawings exemplary constructions of the invention; however, the invention is not  
3 limited to the specific methods and instrumentalities disclosed. In the drawings:

4 FIG. 1 is a block diagram of an exemplary computing environment in  
5 which aspects of the invention may be implemented;

6 FIG. 2 is a block diagram of an exemplary data migration environment  
7 having a computing device communicatively connected to one or more media drives;

8 FIG. 3 is a block diagram showing an exemplary organization for data  
9 objects stored on a medium;

10 FIG. 4A is a block diagram showing a plurality of queues having  
11 requests to recall data objects from media in accordance with aspects of the invention;

12 FIG. 4B is a block diagram of an exemplary queue header data structure;

13 FIG. 4C is a block diagram of an exemplary queue item data structure;

14 FIG. 5 is a flow diagram showing the process by which a recall request  
15 is queued;

16 FIG. 6 is a flow diagram showing an exemplary process for a recall  
17 worker thread;

18 FIG. 7 is a flow diagram showing an exemplary queue activation  
19 technique;

20 FIG. 8A is a diagram showing a queue having items to be recalled whose  
21 offsets form two monotonically increasing sequences;

22 FIGS. 8B and 8C are diagrams showing the insertion of a new item into  
23 the second monotonically increasing sequence of the queue shown in FIG. 8A;

24 FIGS. 8D and 8E are diagrams showing the insertion of a new item into  
25 the first monotonically increasing sequence of the queue shown in FIG. 8A;

26 FIG. 9 is a flow diagram showing an exemplary queue optimization  
27 technique;

1           FIG. 10 is a flow diagram showing an exemplary lock acquisition  
2 process; and

3           FIG. 11 is a block diagram of an exemplary file management  
4 environment in which the invention may be embodied.

5  
6       Overview

7           Many computer systems include a hard disk, or other long-term storage  
8 device, as a primary means for long-term storage of files or other data. When the disk  
9 becomes full, it may be necessary to "migrate" certain data to a backup medium such  
10 as a tape or optical disk. When migrated data is needed by the computer system or its  
11 user, it is necessary to recall that data from the backup medium. The present invention  
12 provides an efficient technique for the recall of data from backup media.

13  
14       Exemplary Computing Environment

15           FIG. 1 illustrates an example of a suitable computing system  
16 environment 100 in which the invention may be implemented. The computing system  
17 environment 100 is only one example of a suitable computing environment and is not  
18 intended to suggest any limitation as to the scope of use or functionality of the  
19 invention. Neither should the computing environment 100 be interpreted as having any  
20 dependency or requirement relating to any one or combination of components illustrated  
21 in the exemplary operating environment 100.

22           The invention is operational with numerous other general purpose or  
23 special purpose computing system environments or configurations. Examples of well  
24 known computing systems, environments, and/or configurations that may be suitable  
25 for use with the invention include, but are not limited to, personal computers, server  
26 computers, hand-held or laptop devices, multiprocessor systems, microprocessor-based  
27 systems, set top boxes, programmable consumer electronics, network PCs,

Computer 110 typically includes a variety of computer readable media. Computer readable media can be any available media that can be accessed by computer 110 and includes both volatile and nonvolatile media, removable and non-removable media. By way of example, and not limitation, computer readable media may comprise computer storage media and communication media. Computer storage media includes both volatile and nonvolatile, removable and non-removable media implemented in any

1 method or technology for storage of information such as computer readable  
2 instructions, data structures, program modules or other data. Computer storage media  
3 includes, but is not limited to, RAM, ROM, EEPROM, flash memory or other memory  
4 technology, CDROM, digital versatile disks (DVD) or other optical disk storage,  
5 magnetic cassettes, magnetic tape, magnetic disk storage or other magnetic storage  
6 devices, or any other medium which can be used to store the desired information and  
7 which can accessed by computer 110. Communication media typically embodies  
8 computer readable instructions, data structures, program modules or other data in a  
9 modulated data signal such as a carrier wave or other transport mechanism and includes  
10 any information delivery media. The term "modulated data signal" means a signal that  
11 has one or more of its characteristics set or changed in such a manner as to encode  
12 information in the signal. By way of example, and not limitation, communication media  
13 includes wired media such as a wired network or direct-wired connection, and wireless  
14 media such as acoustic, RF, infrared and other wireless media. Combinations of any of  
15 the above should also be included within the scope of computer readable media.

16           The system memory 130 includes computer storage media in the form of  
17 volatile and/or nonvolatile memory such as read only memory (ROM) 131 and random  
18 access memory (RAM) 132. A basic input/output system 133 (BIOS), containing the  
19 basic routines that help to transfer information between elements within computer 110,  
20 such as during start-up, is typically stored in ROM 131. RAM 132 typically contains  
21 data and/or program modules that are immediately accessible to and/or presently being  
22 operated on by processing unit 120. By way of example, and not limitation, FIG. 1  
23 illustrates operating system 134, application programs 135, other program modules  
24 136, and program data 137.

25           The computer 110 may also include other removable/non-removable,  
26 volatile/nonvolatile computer storage media. By way of example only, FIG. 1 illustrates  
27 a hard disk drive 140 that reads from or writes to non-removable, nonvolatile magnetic  
28 media, a magnetic disk drive 151 that reads from or writes to a removable, nonvolatile

1 magnetic disk 152, an optical disk drive 155 that reads from or writes to a removable,  
2 nonvolatile optical disk 156, such as a CD ROM or other optical media, and a  
3 sequential media drive 157 that reads from or write to a removable, nonvolatile  
4 sequential medium 158, such as a magnetic tape cassette or reel-to-reel tape. Other  
5 removable/non-removable, volatile/nonvolatile computer storage media that can be used  
6 in the exemplary operating environment include, but are not limited to, flash memory  
7 cards, digital versatile disks, digital video tape, solid state RAM, solid state ROM, and  
8 the like. The hard disk drive 141 is typically connected to the system bus 121 through  
9 an non-removable memory interface such as interface 140, and magnetic disk drive 151  
10 and optical disk drive 155 are typically connected to the system bus 121 by a removable  
11 memory interface, such as interface 150.

12           The drives and their associated computer storage media discussed above  
13 and illustrated in FIG. 1, provide storage of computer readable instructions, data  
14 structures, program modules and other data for the computer 110. In FIG. 1, for  
15 example, hard disk drive 141 is illustrated as storing operating system 144, application  
16 programs 145, other program modules 146, and program data 147. Note that these  
17 components can either be the same as or different from operating system 134,  
18 application programs 135, other program modules 136, and program data 137.  
19 Operating system 144, application programs 145, other program modules 146, and  
20 program data 147 are given different numbers here to illustrate that, at a minimum,  
21 they are different copies. A user may enter commands and information into the  
22 computer 20 through input devices such as a keyboard 162 and pointing device 161,  
23 commonly referred to as a mouse, trackball or touch pad. Other input devices (not  
24 shown) may include a microphone, joystick, game pad, satellite dish, scanner, or the  
25 like. These and other input devices are often connected to the processing unit 120  
26 through a user input interface 160 that is coupled to the system bus, but may be  
27 connected by other interface and bus structures, such as a parallel port, game port or a  
28 universal serial bus (USB). A monitor 191 or other type of display device is also



1 connected to the system bus 121 via an interface, such as a video interface 190. In  
2 addition to the monitor, computers may also include other peripheral output devices  
3 such as speakers 197 and printer 196, which may be connected through an output  
4 peripheral interface 190.

5           The computer 110 may operate in a networked environment using logical  
6 connections to one or more remote computers, such as a remote computer 180. The  
7 remote computer 180 may be a personal computer, a server, a router, a network PC, a  
8 peer device or other common network node, and typically includes many or all of the  
9 elements described above relative to the computer 110, although only a memory storage  
10 device 181 has been illustrated in FIG. 1. The logical connections depicted in FIG. 1  
11 include a local area network (LAN) 171 and a wide area network (WAN) 173, but may  
12 also include other networks. Such networking environments are commonplace in  
13 offices, enterprise-wide computer networks, intranets and the Internet.

14           When used in a LAN networking environment, the computer 110 is  
15 connected to the LAN 171 through a network interface or adapter 170. When used in a  
16 WAN networking environment, the computer 110 typically includes a modem 172 or  
17 other means for establishing communications over the WAN 173, such as the Internet.  
18 The modem 172, which may be internal or external, may be connected to the system  
19 bus 121 via the user input interface 160, or other appropriate mechanism. In a  
20 networked environment, program modules depicted relative to the computer 110, or  
21 portions thereof, may be stored in the remote memory storage device. By way of  
22 example, and not limitation, FIG. 1 illustrates remote application programs 185 as  
23 residing on memory device 181. It will be appreciated that the network connections  
24 shown are exemplary and other means of establishing a communications link between  
25 the computers may be used.

26

## Exemplary Data Migration Environment

FIG. 2 shows an exemplary environment in which migrated data may be stored on backup media. Computer 110 may be communicatively connected to one or more media drives. FIG. 2 shows computer 110 connected to a plurality of media drives 157, which are capable of reading data from media 158 and communicating that data back to computer 110. FIG. 2 depicts media drives 157 and media 158 as tape drives and magnetic cassette tapes, respectively. Tape media, however, is merely exemplary and is not limiting of the invention. The environment shown in FIG. 2 could, as an alternative example, include magnetic disk drives 151 for use with magnetic disks 152 (shown in FIG. 1), optical disk drives 155 for use with optical disks 156 (also shown in FIG. 1), or any other type of media-reading device for use with appropriate data storage media.

Media drives 157 are preferably configured to work with "removable" media, such that a given medium 158 can be mounted or dismounted from drive 157. A media drive 157 may, at any given time, have a particular medium 158 mounted on it, or it may be "empty" (i.e., no medium 158 presently mounted on the drive). By way of example, FIG. 2 shows one media drive 157 (the leftmost media drive 157 in FIG. 2) having a medium 158 mounted thereon, and three other media drives 157 have no media mounted thereon.

Media drives 157 may be included within, or may be associated with, a "juke box" 202. Juke box 202 stores a library 204 of media 158 in a "near-line" position such that media 158 in library 204 can be mounted and dismounted by machine without human intervention. Juke box 202 includes robotic arm 206, which mounts media onto (and dismounts media from) media drives 157. While juke box 202 is a convenient device for storing, mounting, and dismounting media 158, juke box 202 is optional, and the invention applies to any environment having one or more media drives 157, whether or not media drives 157 or media 158 are associated with juke box 202.

Media 158 are generally used to store "migrated" data. In a typical example, media 158 are used to store files that have been "migrated" (i.e., moved off of) a primary storage device (such as hard disk 141 shown in FIG. 1) in order to free up space for new data on the primary storage device. For example, a storage management program on computer 110 (which may, for example, be part of operating system 134) may target files on hard disk 141 that have not been accessed for some predetermined period of time (e.g., six months). The targeted files may be migrated to media 158 by deleting the files from hard disk 141 and copying them to one or more media 158. A "stub" that identifies the new location of each file may be retained on hard disk 141 so that the file can be located later. While old or rarely-used files are a common type of data that are stored on media 158, it should be understood that such files are not limiting of the invention. Rather, the invention may be used to recall any type of data that may be stored on media 158, regardless of whether the data is organized into files, and regardless of the reasons for which the data has been stored on media 158.

#### Exemplary Organization of Data Objects on a Medium

FIG. 3 depicts an example organization of data on a medium 302. Medium 302 could be any type of medium, such as tape 158, optical disk 156, or magnetic disk 152 (all shown in FIG. 1). In this example, medium 302 stores files 304, although it should be appreciated that a file is merely an exemplary type of data object, and medium 302 could store any type of data objects. As noted above, files 304 may, for example, be files that were once resident on a hard disk (e.g., hard disk 141 shown in FIG. 1) but that have been "migrated" to medium 302 in order to free up space on the hard disk. It may be convenient to group files 304 into "bags" 306 for storage on medium 302. Essentially, a "bag" is an arbitrary grouping of files that are stored together. As one example, the files 304 in a given bag 306 may represent all files from a particular volume of hard disk storage that were migrated at a particular time. Thus,

1 if the migration process is performed nightly, then each bag 306 may represent all of  
2 the files from a given volume that were collected in a given night (e.g., all files from  
3 volume C: that were migrated on 1/1/2000). It should be understood that the  
4 organization of data objects into bags 306 is merely for convenience, and data objects  
5 could, alternatively, be stored individually on medium 302 without grouping those data  
6 objects into bags 306.

7           The storage locations in medium 302 may have a one-dimensional linear  
8 order. That is, every storage location on medium 302 may be described by a single-  
9 number offset from a starting position. This order may be implemented physically, as  
10 in the case of sequential media such as tapes where the offset represents essentially the  
11 linear distance from the start of the tape. Alternatively, the order may be implemented  
12 "logically," as in the case of a magnetic disk in which an arbitrary "start" point is  
13 defined on each track, and in which the tracks are traversed in a predefined sequence,  
14 thereby defining an order for the entire set of storage locations on the disk. Thus, as  
15 shown in FIG. 3, every file 304 is shown as starting at a location that is described by its  
16 offset from the "beginning" of the medium, whether that "beginning" is a physical or  
17 logical beginning. In FIG. 3, the beginning of the medium is defined as offset zero.  
18 File AA begins at offset 1010, file AB begins at offset 1050, and so on. Moreover,  
19 each bag 306 can be described as beginning at a particular offset into medium 302: bag  
20 A begins at offset 1000, bag B begins at offset 2000, and so on. Thus, if the starting  
21 offset for each bag is known, a file's offset into medium 302 can be described by the  
22 identity of the bag 306 in which it is located and the file's offset into the bag. For  
23 example, file AA is in bag A, and has offset 10 into bag A. Thus, adding 1000 (the  
24 bag's offset into the medium) together with 10 (the file's offset into the bag), produces  
25 the absolute offset (1010) of file AA into the medium.

26



1 linked list, or another data structure that stores items 406 in an identifiable sequence  
2 within a queue. As further discussed below in connection with FIG. 8A-9, certain  
3 techniques may be used to optimize the order in which items are placed in the queue.

4 Each queue has a status 404 associated with it, which can be either active  
5 or non-active. In the example of FIG. 4A, queue B is active and queues A and Z are  
6 non-active. Queue B's active status indicates that medium B is presently mounted on a  
7 drive and is being used to recall the data objects indicated by the items 406 on queue B.  
8 The inactive status of queues A and Z indicates that media A and Z are not presently  
9 being used to recall the data objects indicated by the items 406 on those queues. For  
10 example, queues A and Z may be inactive because there is no available media drive  
11 onto which to mount media A and/or Z. As another example, media A and/or Z may be  
12 presently mounted elsewhere and being used for a different purpose, or may be stored  
13 off-line such that they are not available for mounting.

14 Each queue 402 may be represented by a queue header. FIG. 4B shows  
15 an exemplary queue header data structure 420. The exemplary data structure includes a  
16 media identifier 422, a timestamp 424, a list pointer 426, and a lock 428.

17 Media identifier 422 identifies the particular medium to which the queue  
18 402 corresponds. With reference to the example of FIG. 4A, media identifier 422 could  
19 be "A", thereby indicating that the requests in the queue 402 defined by queue header  
20 420 are requests for data objects located on medium A. In general, media identifier 422  
21 is usually a number or character string that identifies the medium.

22 Timestamp 424 represents the time at which the queue 402 was created.  
23 As further discussed below, when it is necessary to choose a non-active queue 402 to  
24 activate, it may be advantageous to activate the oldest non-active queue 402 first.  
25 Timestamp 424 allows the oldest queue 402 to be identified.

26 Lock 428 is a read/write lock on queue 402. Each queue 402 has a lock  
27 428. The per-queue lock is used for locking the queue while searching for the position  
28 to insert an item into the queue. Many items may arrive to the migration engine at the

1 same time. However, the search for position and insertion of an item is protected by a  
2 lock in order to ensure that the item is inserted at the correct place in the queue. Lock  
3 428 ensures that items are inserted in a given queue one at a time.

4 Item pointer 426 contains the address of the first item 406 in the queue.  
5 As noted above, items in the queue may be represented as a linked list. FIG. 4C shows  
6 an exemplary structure 440 for representing items 406 as a linked list. Linked list  
7 element structure 440 includes data 442 indicating the location on a medium 158 of a  
8 requested data object. Structure 440 also includes a link 444 to the next item 406 on the  
9 list. Link 444 preferably takes the form of a pointer to another linked list element  
10 structure 440. When items 406 are represented in this manner, then item pointer 426 in  
11 FIG. 4B is simply a pointer to the first linked list element structure 440 in the queue  
12 402. However, it should be understood that items 406 could be represented by another  
13 type of data structure, such as an array. In such a case, item pointer 426 may contain  
14 the address of the beginning of the array.

#### 15 16 Recall Technique: Placing a Request on One of a Plurality of Queues

17 With reference to FIG. 5, the procedure will now be described by which  
18 recall requests are queued for processing. First, a request to recall a data object (e.g.,  
19 file 304) from media 302 is received (step 502). The request may take any format that  
20 identifies the requested data object. For example, the request may identify a particular  
21 data object by its bag ID, and its offset into the bag. However, the request may take  
22 any form that identifies the data object requested. For example, the request may specify  
23 the particular medium 302 (e.g., medium A in FIG. 4A) on which the requested data  
24 object is located and its absolute offset into that medium 158. The request received at  
25 step 502 may identify the requested data object in any manner without departing from  
26 the spirit and scope of the invention.

27 At step 504, the absolute location of the data object (i.e., its media ID  
28 and offset into the medium) is ascertained based on the identifying information in the

At step 508, the set of existing queues 402 is examined, and it is determined whether a queue exists for the medium on which the requested data object is located. This action may be performed by examining the media ID field 422 in the queue headers 420 of all of the existing queues 402. If it is determined that such a queue exists, then the request received at step 502 is added to the appropriate queue (step 526), the lock is released (step 528), and the process terminates.



1           On the other hand, if it is determined at step 508 that no such queue  
2 exists, then a new queue is created at step 510. The creation of a new queue may be  
3 performed by creating a new queue header data structure 420, noting the time in  
4 timestamp field 424, and assigning the appropriate media ID to media ID field 422. The  
5 request received at step 502 is then added to the new queue at step 510.

6           When step 508 results in a determination that a new queue needs to be  
7 created, it must be determined whether a new queue should be activated. This  
8 determination is made at step 514 by determining whether the number of active queues  
9 (which may be maintained in a global counter) is less than a global value called the  
10 "concurrency." The concurrency is the maximum number of media that can be  
11 simultaneously used for recall, and thus is also the maximum number of queues that can  
12 be active at a given point in time. The concurrency is partly bounded by the physical  
13 parameters of the environment. Referring for the moment back to FIG. 2, it will be  
14 recalled that a system may have only a finite number of media drives 157 available to  
15 it. Since only one medium may be mounted on a media drive 157 at a given time, the  
16 number of media drives accessible to the system performing the recall is an upper limit  
17 on the concurrency. However, the concurrency may be a lower number, and may be  
18 user-specified. For example, the system administrator may decide that one drive should  
19 always be kept free for non-recall use, in which case he may set the concurrency at one  
20 less than the number of physical drives. For example, although FIG. 2 depicts  
21 computer 110 connected to four media drives 157, it may be the case that the system  
22 operator has designated only three of those drives for recall operations, and thus the  
23 concurrency in such a case would be three.

24           It should be observed that the number of queues that are active at a given  
25 time is usually not greater than the concurrency since (as more particularly discussed  
26 below) a new queue will only be activated when the number of presently active queues  
27 is less than the concurrency. However, it is possible in some cases for the number of  
28 active queues to exceed the concurrency when the concurrency is adjusted dynamically

1 (e.g., by a system administrator) during the operation of the system. For example, if  
2 the concurrency is set to three, and three media are presently mounted on three  
3 different drives for recall of data objects, at some point during the recall operation a  
4 system administrator may lower the concurrency to two. In such a case, the presently  
5 active queues may be allowed to finish, but, until one queue completes processing, the  
6 number of active queues will be greater than the concurrency.

7           Returning now to FIG. 5, if it is determined at step 514 that the number  
8 of active queues is not less than the concurrency, then this means that the maximum  
9 number of queues is already active so no other queues can be activated. In this case, the  
10 global lock is released (step 516), and the process ends.

11           On the other hand, if it is determined at step 514 that the number of  
12 active queues less than the concurrency, then one of the inactive queues (i.e., the queue  
13 that was newly-created at step 510, or another inactive queue) can be activated. In  
14 order to activate a new queue, a counter that represents the number of presently active  
15 queues is incremented at step 518. Then, the global lock is released at step 520. At step  
16 522, the medium corresponding to the queue is mounted on one of media drives 157.  
17 At step 524, a new "recall worker thread" (described in FIG. 6) is spawned in order to  
18 process the newly activated queue.

19  
20 Recall Technique: Exemplary Recall Worker Thread

21           As previously noted, one advantage of the invention is the ability to  
22 concurrently recall data objects from different media when plural media drives are  
23 available. Concurrent recall is performed by establishing a separate "worker thread"  
24 for each active queue 402. The "worker thread" is established at step 524 in FIG. 5, or  
25 at step 716 in FIG. 7 discussed below. FIG. 6 shows an exemplary process for a "recall  
26 worker thread" established at one of the aforesaid steps.

27           Referring now to FIG. 6, at the start of the worker thread process the  
28 first item on the queue (e.g., item 406 in FIG. 4A) is popped from the top of the queue

1 (step 602). The item is evaluated to determine the location on the medium of the next  
2 data object to be recalled. When the top item is popped from the queue, the next item  
3 on the queue becomes the top item.

4 At step 604, the data object identified by the popped item is read from  
5 the medium. At step 606, the data object that was read at step 604 is written to a target  
6 location. For example, if the data objects stored on the medium are migrated files, then  
7 the data read from the medium at step 604 may be written back onto a hard disk (e.g.,  
8 hard disk 141 in FIG. 1) into space that is designated for the storage of the migrated  
9 file.

10 At step 608, it is determined whether the queue being processed by the  
11 worker thread is empty – that is, whether the last item has been popped off the queue.  
12 If the queue is not empty, the process returns to step 602, where the process heretofore  
13 is repeated for the next item on the queue. If the queue is empty, then the global lock is  
14 acquired at step 610, and the process continues to step 612.

15 At step 612, it is again determined whether the queue is empty. Although  
16 the queue has previously been tested at step 608 to determine whether it is empty, it is  
17 preferable to perform that test again due to the small chance that another thread may  
18 have placed another item on the queue after the test was last performed but before the  
19 lock was acquired – i.e., between steps 608 and 610. For example, a concurrently  
20 executing thread performing the queuing process depicted in FIG. 5 might have placed  
21 another item on the queue after step 608 is performed but before step 610 is performed.  
22 If the determination at step 612 is that the queue is not empty (i.e., has become non-  
23 empty since step 608), then the lock is released at step 614 and the process returns to  
24 step 602 to process the items that have been placed on the queue.

25 If step 612 results in a determination that the queue is still empty, then  
26 the counter that maintains the number of active queues is decremented (step 616). The  
27 counter decremented at step 616 is the same counter that is incremented at step 518 in  
28 FIG. 5. After decrementing the counter, the queue is destroyed (step 618). The exact

1 manner in which the queue is destroyed depends on the manner in which the queue is  
2 represented. For example, if the queue is represented by a queue header 420 (shown in  
3 FIG. 4B), then destruction of the queue may be accomplished by deallocating the  
4 memory reserved for queue header 420.

5 Next, at step 620 a procedure is started to activate a new queue. An  
6 exemplary queue activation procedure is described below in connection with FIG. 7. It  
7 should be noted that step 620 does not necessarily result in the activation of a new  
8 queue, since there is a possibility that (1) the concurrency has been dynamically  
9 lowered since the last queue activation and thus the current number of active queues is  
10 equal to (or greater than) the present concurrency, or (2) there are no non-active queues  
11 waiting to be activated. Tests for these conditions are more particularly described below  
12 in connection with FIG. 7. Thus, step 620 results in an attempt to activate a new queue,  
13 and, depending upon conditions, may actually result in the activation of a new queue.

14 Following step 620, the global lock is released (step 622). Subsequently,  
15 the "recall worker thread" is destroyed at step 624, although it will be noted that a new  
16 "recall worker thread" may have spawned to replace it at step 620.

#### 17 18 Exemplary Queue Activation Technique

19 FIG. 7 shows an exemplary process for activating a queue. The process  
20 shown in FIG. 7 is a essentially preferred method of performing step 620 (shown in  
21 FIG. 6).

22 At step 702, it is determined whether there are any non-active queues. If  
23 there are no non-active queues, then the process terminates. If there are non-active  
24 queues, then the process proceeds to step 704.

25 At step 704, a queue is selected from among the existing non-active  
26 queues. In the embodiment of the invention in which each queue has a timestamp as  
27 shown in FIG. 4B, it may be useful to select the queue having the oldest timestamp.  
28 However, it will be understood that such a selection paradigm is merely exemplary, and

1 a queue may be selected at step 704 in any manner without departing from the spirit  
2 and scope of the invention. Once a queue is selected for activation, the global lock is  
3 acquired at step 706.

4 At step 708, it is determined whether the number of active queues is less  
5 than the concurrency. If the number of active queues is not less than the concurrency,  
6 then the process continues to step 718 to release the lock and the queue activation  
7 process ends without activating a queue. If the number of active queues is less than the  
8 concurrency, then the counter that maintains the number of active queues is  
9 incremented (step 710). The counter incremented at step 710 is the same counter  
10 referred to in step 518 (FIG. 5) and step 616 (FIG. 6). After the counter is  
11 incremented, the lock is released (step 712).

12 At step 714, the medium 158 corresponding to the queue 402 selected at  
13 step 704 is mounted on one of media drives 157. The mounting may be performed in  
14 any manner appropriate for the environment in which the recall takes place. For  
15 example, if juke box 202 is present (as shown in FIG. 2), then mounting the selected  
16 medium may be performed simply by instructing juke box 202 to mount the selected  
17 medium using robotic arm 206. Alternatively, if no juke box 202 is present, or if the  
18 selected medium is stored outside of the library 204 that is accessible to robotic arm  
19 206, then mounting the selected medium may be performed by prompting a human  
20 operator to mount the medium.

21 At step 716, a new thread is started to process the activated queue. The  
22 started thread is a "recall worker thread," as depicted in FIG. 6. After the new "recall  
23 worker thread" is started, the queue activation process terminates.

#### 24 25 Queue Optimization Technique

26 Referring back to step 512 in FIG. 5 wherein requests for items are  
27 placed on a queue 402, the requests may be placed on the queue in any order. For  
28 example, requests may be placed on the end of the queue in the order in which they are

1 received without regard to the location on the medium of the requested data object.  
2 However, queuing of items in this manner has the disadvantage that it may cause the  
3 medium to shuttle back and forth across the reading head as data objects are accessed in  
4 an order that does not take into account where the data objects are located on the  
5 medium. The problem is particularly acute in the case of tape media, where such  
6 shuttling not only increases the amount of time required for recall operations but also  
7 increases wear on the tape. One way to improve performance is to place items on a  
8 queue such that the requests on the queue form up to two sequences, where each  
9 sequence increases monotonically with respect to the offsets of the requested data  
10 objects.

11 FIG. 8A shows an example of items on a queue 402 organized into two  
12 monotonically increasing sequences. The queue represents requests for data objects on a  
13 particular medium, where there are six different requested data objects located at offsets  
14 50, 75, 90, 100, 200, and 300. These requests are organized into a first sequence 802  
15 and a second sequence 804. It will be observed that the items within a given sequence  
16 are always increasing with respect to their offsets. That is, for any item within a given  
17 sequence, the next item within that same sequence always has a larger offset. (It should  
18 be noted that a subset of each sequence (e.g., the sequence 100, 200) is also a  
19 “monotonically increasing sequence.” However, as used herein the term  
20 “monotonically increasing sequence” refers to the largest such sequence – i.e., a  
21 sequence such that adding any of the adjacent items to the sequence would cause the  
22 sequence not to be monotonically increasing. Such a “largest” sequence can be  
23 described as a “maximal” monotonically increasing sequence.)

24 A purpose of using two monotonically increasing sequences is to address  
25 the situation in which a request for a data object is dynamically queued during the  
26 reading of the relevant medium, and where the reading head has already passed the  
27 newly-queued object’s location on the medium. Essentially, the second sequence 804  
28 represents a second “pass” through the medium that will be made to recall such

1 dynamically-queued requests. Thus, instead of shuttling back and forth across the  
2 medium (which might happen if requests were placed on the queue in no particular  
3 order) the medium is traversed once in a single direction to process the requests in  
4 sequence 802, and then a second time to process the requests in sequence 804.

5 FIG. 8B shows an example of a circumstance in which a request is  
6 inserted into second sequence 804. Queue 402 depicted in FIG. 8B is the queue that  
7 corresponds to a given medium 158, and reading head 850 is positioned at offset 300  
8 along medium 158. At that point in time, a new request is queued, and the request is  
9 for a data object located at offset 215 along medium 158. Because reading head 850 has  
10 already passed location 215 while making its pass through first sequence 802, the item  
11 is not queued in first sequence 802 but rather is placed in second sequence 804, as  
12 shown in FIG. 8C. It should be observed that the new item is positioned in second  
13 sequence 804 such that the increasing nature of the offsets in each sequence is  
14 maintained (in FIG. 8C, the new request having offset 215 is placed after the request  
15 for a data object located at offset 90).

16 FIG. 8D shows an example of a circumstance in which a request is  
17 inserted into first sequence 802. In FIG. 8D, the new item to be queued is for a data  
18 object located at offset 215, as in FIG. 8B. Unlike FIG. 8B, however, FIG. 8D shows  
19 reading head 850 positioned at offset 200. Thus, the new item having offset 215 can be  
20 queued in first sequence 802 (as shown in FIG. 8E), because reading head 850 still has  
21 yet to pass over location 215 while making its pass through first sequence 802.

22 FIG. 9 shows the process of inserting a request into a queue that is  
23 organized into two monotonically increasing sequences, as described above. At step  
24 901, a determination is made as to whether the offset of the requested data object into  
25 the medium is greater than the current offset of the reading head. When the medium is  
26 mounted, the current offset is equal to the actual position of the reading head along the  
27 medium. When the medium is not mounted, the current offset is set equal to zero.  
28 While a non-mounted medium is not technically positioned at any offset relative to the

1 reading head, it may be viewed as being positioned at a zero offset. Since reading will  
2 commence at the starting position when such a medium is mounted, the entire content  
3 of the medium effectively lies ahead of the reading head, because when the medium is  
4 mounted reading will commence at the starting position. (A possible exception is  
5 cartridge tape media which can be dismounted while the tape is wound to any position,  
6 in which case the last known offset can be stored in a memory location.)

7           If step 901 results in a determination that the requested data object's  
8 offset into the medium is greater than the current reading head offset, then the request  
9 is inserted into the first sequence (step 902). The request is inserted into the sequence in  
10 such a position that the monotonic increasing nature of the sequence is preserved - i.e.,  
11 after an item having a lower offset but before an item having a higher offset. If it is  
12 determined at step 901 that the offset of the requested item is less or equal to than the  
13 current reading head offset, then the item is inserted into the second sequence (step  
14 903). After inserting the request into either sequence, the process terminates.

15           It should be observed that when a queue 402 is non-active (i.e., when  
16 reading of its corresponding medium 158 has not yet begun), all newly queued requests  
17 will be placed in a single sequence. This is so because there is no possibility that  
18 reading head 850 has passed the location of the newly-requested data object because  
19 reading of the medium has not even begun (except, as noted above, where a non-  
20 mounted cartridge tape is wound to a non-zero offset). In this sense, the situation in  
21 which items are added to a non-active queue is essentially equivalent to the situation in  
22 which reading head 850 is located at the beginning of the medium (i.e., at offset zero).  
23 Additionally, even if reading head 850 is not located at offset zero, if all incoming  
24 requests are located after the current position of reading head 850 then there will only  
25 be one sequence in the queue.

26           Additionally, it should be observed that the two monotonically increasing  
27 sequences in each queue are actually a series of monotonically increasing sequences that  
28 are being dynamically created and exhausted. For example, as a queue is being



1 processed, its second sequence may begin to grow as new items, whose locations are  
2 behind the read head, are dynamically placed on the queue. Once the first sequence is  
3 exhausted, the first sequence no longer exists, because all items in the first sequence  
4 have been popped off the queue. At this point, the “second” sequence in effect becomes  
5 the first (and only) sequence in the queue. However, once processing of this sequence  
6 begins, new items may be dynamically queued that are behind the read head. These  
7 items are placed in a “new” second sequence. While this new sequence is the second  
8 sequence of items that are presently on the queue, it is historically the third sequence of  
9 items that have been created for the queue. However, since all of the items in the first  
10 historical sequence were popped off the queue during the first pass through the  
11 medium, the new sequence is the second sequence among those sequences that presently  
12 exist on the queue. Thus, the characterization of the queue as having two sequences  
13 refers to the state of the queue at any given point in time, rather than an historical  
14 description of what sequences have ever existed on the queue.

15

#### 16 The Global Lock

17 It will be observed with reference to FIGS. 5-7 that a global lock is used  
18 (e.g., at steps 506, 520, 528, 610, 614, 622, 706, 712, and 718) in order to protect  
19 certain operations. Because any of the processes depicted in FIGS. 5-7 could be  
20 operating concurrently on different threads, the global lock is used to prevent errors  
21 that could result if two or more concurrent threads performed certain operations at the  
22 same time. The use of a lock guarantees that, even in a concurrent execution  
23 environment, certain classes of operations (i.e., those surrounded by locking and  
24 unlocking operations) will not be performed concurrently by two different threads.

25 FIG. 10 shows how the various “acquire lock” steps are performed (i.e.,  
26 steps 506 in FIG. 5, step 610 in FIG. 6, and step 706 in FIG. 7). FIG. 10 is essentially  
27 a detailed description of what happens at steps 506, 610, and 706. First, a  
28 determination is made as to whether the lock is available (step 1002). The lock is

1 unavailable if another thread has acquired it but has not yet released it; otherwise, it is  
2 unavailable. If it is determined at step 1002 that the lock is available, then the thread  
3 proceeds at step 1006 to perform whatever is its next step after lock acquisition. For  
4 example, in the queuing procedure of FIG. 5, lock acquisition is performed at step 506,  
5 so if the lock is available then the process is permitted to proceed to the step following  
6 step 506 (i.e., step 508). On the other hand, if it is determined at step 1002 that the  
7 lock is unavailable, then execution switches to another thread for some amount of time  
8 (step 1004). Eventually, control returns to the thread that is waiting for the lock, where  
9 step 1002 is performed again in order to determine whether the lock has become  
10 available. The cycle of testing for the lock's availability (at step 1002) and executing  
11 another thread for some amount of time (at step 1004) is repeated until the lock  
12 eventually becomes available.

13 In the present invention, the use of the global lock is designed to  
14 prevent: (1) the destruction of a queue by a first thread (at step 618) while a second  
15 thread is queuing a request on that queue (at step 526); and (2) the simultaneous  
16 activation of two different queues (e.g., by steps 518-524 or steps 710-716) that could  
17 otherwise result in an active queue count that exceeds the concurrency.

18 For example, suppose that a first thread is executing the queuing  
19 procedure shown in FIG. 5, and a second thread is executing the "recall worker thread"  
20 shown in FIG. 6. If both threads could proceed concurrently without regard to each  
21 other (i.e., if the locking and unlocking steps in those procedures were not performed),  
22 then the situation might arise in which the first thread identifies a particular queue for  
23 queuing a new request (at step 508), and then the execution context switches to the  
24 second thread, which destroys that same queue at step 618. When control switches back  
25 to the first thread to place the request on the queue (at step 526), the queue would  
26 already have been destroyed, resulting in an error. However, the use of a global lock  
27 prevents the possibility of this error, because the second thread can never execute the  
28 queue destruction step (step 618), or even perform the final test for queue emptiness

007330 034 E34

1 (step 612), while the first thread is performing steps 508 and 526, since those steps are  
2 protected by the same lock.

3           As another example, suppose that a first thread is executing the queuing  
4 procedure shown in FIG. 5, and a second thread is executing the queue activation  
5 procedure shown in FIG. 7. Suppose that the first thread has just created a new queue  
6 (at step 510) and added a new request to the new queue (at step 514). Moreover,  
7 suppose that the number of active queues is one less than the concurrency (i.e., there  
8 are sufficient drives available to active one additional queue). The first thread then  
9 proceeds to test whether the number of active queues is less than the concurrency (at  
10 step 514), and concludes that it can activate a new queue. If both the first and second  
11 threads could proceed concurrently without regard to each other, then it is possible that  
12 immediately after the first thread tests the value of the active queue counter (at step  
13 514), the execution context would switch to the second thread, which would also test  
14 the value of the queue counter (at step 708). Since the first thread has not yet updated  
15 the active queue counter (at step 518), the second thread will also conclude that it can  
16 activate a new queue. In this situation, two threads will each proceed to activate a new  
17 queue, even though there are only sufficient resources to active one queue. This  
18 situation is prevented, however, by the use of the lock: since the first thread does not  
19 release the lock (step 520) until after it has incremented the active queue counter (step  
20 518), the second thread will not reach the step of examining the active queue counter  
21 (at step 708) until the first thread has appropriately updated the counter (at step 518)  
22 and released the lock (at step 520).

23           Thus, the use of a lock, while not mandatory, is a useful way of  
24 preventing certain types of errors that could occur if two or more threads performed  
25 certain operations at the same time. However, other methods of synchronizing  
26 operations among concurrently executing threads, and such other methods may be used  
27 without departing from the spirit and scope of the invention.

13

## 15

25

1 monotonically increasing sequences (as shown in FIGS. 8A-9), and recall was  
2 performed concurrently from both available drives (i.e., the concurrency was set to 2).  
3 When the same set of 200 recall requests, issued in the same order as in the first test,  
4 were processed in this manner, processing of the recalls took only 5 minutes and 7  
5 seconds – an approximately nine-fold increase in throughput.

### 6 7 Exemplary File Management Environment

8 While the present invention may be used in any context to recall data  
9 stored on media, a particularly useful environment in which the invention may be  
10 embodied is a file management environment which manages files, some of which have  
11 been “migrated” to backup media. FIG. 11 shows such an environment in which the  
12 invention may be incorporated.

13 Referring now to FIG. 11, a file management environment 1100 (which  
14 may be part of an operating system, such as operating system 134 shown in FIG. 1)  
15 comprises a remote storage filter 1102, a file system manager 1104 (of which NTFS,  
16 depicted in FIG. 11, is an example), and one or more volumes of files 1106. Volumes  
17 1106 may, for example, be plural hard disks or separate partitions of a single hard disk.  
18 File management environment is organized in a hierarchy: requests and instructions are  
19 received by file management system from the user level by way of remote storage filter  
20 1102. Remote storage filter 1102 receives requests for files (such as request 1108 to  
21 open a file called “foo,” depicted in FIG. 11). Remote storage filter 1102 passes the  
22 request to file system 1104, which, in turn, locates the requested file on the appropriate  
23 volume 1106. File system 1104 includes data or logic that is able to identify the  
24 particular volume 1106 on which the requested file is stored. File system 1104 may  
25 then pass the file 1110 retrieved from volume 1106 back to user mode by way of RS  
26 filter 1102. The file is then made available in user mode.

27 File system 1104 may maintain a set of reparse points 1112. A reparse  
28 point is essentially a flag and a storage location for arbitrary data. When a request for a

007283830 "004504"

1 file is received (e.g., a request for the file "foo"), file system 1104 may check to  
2 determine whether a reparse point is set for that file. If there is no reparse point 1112  
3 for the file, file system 1104 locates the file in volumes 1106. If a reparse point 1112 is  
4 set for the file, then file system 1104 indicates this fact back to remote storage filter  
5 1102, along with whatever arbitrary data is associated with the particular reparse point.  
6 The indication that reparse point 1112 has been set serves as a flag to remote storage  
7 filter 1102 indicating that the requested file is not located in volumes 1106. The  
8 arbitrary data associated with reparse point 1112 may be a "stub" that indicates where  
9 the file is stored in remote storage - e.g., a media ID, a bag ID, and an offset.  
10 Specifically, the process of migrating files to remote storage may include the step of  
11 setting a reparse point and storing in the reparse point the remote location of the file.  
12 Remote storage filter 1102 then communicates with remote storage engine 1114, which  
13 receives the information identifying the location of the requested file, retrieves that file  
14 from media using a physical remote storage device (e.g., juke box 202) and provides  
15 that file back to remote storage filter 1102. Remote storage filter then takes the file  
16 received from remote storage engine 1114 and provides it to user mode. The process of  
17 hitting a reparse point 1112 and retrieving a file from remote storage in response  
18 thereto may, in fact, take place "transparently." That is, when a user issues a request  
19 1108 for a file, the components of file management environment 1100 may act together  
20 "behind the scenes" such that the user is unaware of whether the file was received from  
21 volumes 1106 or from a remote storage medium (except, of course, that retrieval from  
22 remote storage may take a longer time).

23           Features of the present invention may be embodied in remote storage  
24 engine 1114. Generically, remote storage engine 1114 is a software component that  
25 contains logic which retrieves requested data from media using a physical remote  
26 storage device. Conventionally, remote storage engine 1114 may contain logic which  
27 queues requests sequentially (in a single queue) and processes the requests in the order  
28 received. In accordance with the invention, however, remote storage engine 1114 may

00783860 "0004504

1 be programmed with logic that retrieves requests according to the technique described  
2 above in connection with FIGS. 4A-9. That is, remote storage engine 1114 may  
3 perform the operations of creating and managing plural queues (one for each medium  
4 on which requested data objects are located), activating and destroying queues, storing  
5 a concurrency value and various counters, organizing the requests on each queue to  
6 form two monotonically increasing sequences, and all of the other techniques depicted  
7 in FIG. 4A-9.

8           When aspects of the invention are embodied in remote storage engine  
9 1114, requested files that have been migrated to media may be retrieved in the  
10 following manner. First a user issues a request to access file, such as "open foo"  
11 request 1108. This request is received by remote storage filter 1102, which passes the  
12 request to file system 1104. If the file "foo" has been migrated to media and is not  
13 located on volumes 1106, then a reparse point 1112 may have been set for the file.  
14 Thus, file system 1104 notifies remote storage filter 1102 that a reparse point 1112 has  
15 been set for the requested file, and also provides to remote storage filter 1102 the data  
16 that is stored with reparse point 1112. As noted above, this data stored with reparse  
17 point 1112 may include a "stub" indicating the file's location on media. Remote storage  
18 filter 1102 then issues a request for the file to remote storage engine 1114. This request  
19 is the request received at step 502 of FIG. 5. Remote storage engine 1114 then queues  
20 the request according to the process of FIG. 5, creating a new queue if necessary. It  
21 should be noted that database 550 (depicted in FIG. 5) may be accessible to remote  
22 storage engine 1114. Database 550 may be used to convert certain type of file  
23 identifying information (i.e., bag ID, and bag offset), into an absolute location for the  
24 file (i.e., a media ID and a media offset). Additionally, inasmuch as plural versions of  
25 the same file may have been migrated on different occasions, the information stored in  
26 database 550 may be used to ascertain which is the most recent version of the file (and  
27 thus the version that should be recalled). Remote storage engine 1114 then processes  
28 the requests by communicating with a physical remote storage device (e.g., a juke box)

0070300 "004604  
"0070300 "004604

It is noted that the foregoing examples have been provided merely for the purpose of explanation and are in no way to be construed as limiting of the present invention. While the invention has been described with reference to various embodiments, it is understood that the words which have been used herein are words of description and illustration, rather than words of limitations. Further, although the invention has been described herein with reference to particular means, materials and embodiments, the invention is not intended to be limited to the particulars disclosed herein; rather, the invention extends to all functionally equivalent structures, methods and uses, such as are within the scope of the appended claims. Those skilled in the art, having the benefit of the teachings of this specification, may effect numerous modifications thereto and changes may be made without departing from the scope and spirit of the invention in its aspects.